

CROSS-REFERENCE TO RELATED APPLICATION

(1) U.S. Patent Application Serial No. ____/____ (Attorney Docket No. AUS920040005US1).

10 BACKGROUND OF THE INVENTION

1. Technical Field:

15 **[0001]** The present invention relates in general to improved high availability cluster management and, in particular to remote cluster management of a high availability system. Still more particularly, the present invention relates to improved remote monitoring and management of multiple high availability systems in an enterprise network.

2. Description of the Related Art:

[0002] For retailers, banks, and other on-line services where load and demand

constantly fluctuate and where handling each customer request is of utmost importance, high availability (HA) systems have been developed to handle mission-critical operations. In general, an HA system is a system designed to eliminate or minimize the loss of service due to either planned or unplanned outages among components of a network system. The key method of providing an HA system is through redundant hardware and software components grouped into a cluster of servers.

[0003] Redundancy is important in an HA system because when a failure occurs in one node of the cluster, the system transfers the processes performed by one node to another. In a two-node HA cluster, for example, one node is typically designate as the primary node and the other node is typically designated as the backup node. In general, the primary node initially runs an application when a cluster is started. In addition, in general, a backup node is designated that will run the application if the primary node fails. The HA cluster system will typically implement a cluster manager process that periodically polls (or checks the heartbeat) of the primary node to determine if it is still active. If a “heartbeat” is not detected, then the cluster manager moves the software process to another server in a cluster.

[0004] An important characteristic of an HA system is the recovery time. In general, the recovery time in a HA system is the time taken for a backup node to take over an application from a failed primary node. Recovery time is particularly important in a sales based HA system because retailers may lose valuable business if a customer is not able to complete transactions quickly. A delay of even 30 seconds for the recovery time diminishes a retailer’s business transactions.

[0005] Another important characteristic of an HA system is to achieve little or no loss of data during failover. In particular, it is important to achieve little or no loss of committed data. For example, it is not advantageous to lose valuable information about a customer order or customer information during failover.

5 [0006] To achieve a short recovery time and little or no loss of data during failure, it is important to initially combine hardware and software in such a manner that an HA system is built. After a HA system is initiated, however, it is important to monitor and adjust the configuration of the HA system to try to improve the efficiency of failovers and correction of other errors.

10 [0007] When configuring hardware and software for HA systems, many developers have developed customized HA software services to control applications in a custom environment which often requires new hardware. These solutions are often expensive and do not take advantage of open source technologies that allow for portability of applications across multiple platforms. Further, expensive server systems are often selected, in hopes that the power
15 available in the server system will automatically increase the efficiency of failovers.

[0008] As an alternative, open source developers continue to expand open source technology with functions that can be configured when implementing HA systems. For example, Linux provides an inexpensive, platform independent operating system. Developers of Linux continue to add functions to the operating system that can be implemented in an open
20 source manner by other developers. Some of these functions, such as “heartbeat” and distributed replicated block device (drbd), are implemented with the Linux operating system to assist in

configuring HA systems.

[0009] While the Linux tools provide a framework for monitoring for failures and configuring the hardware used in HA systems, there is a need for additional monitoring and configuration capability. In particular, there is a need for a method of monitoring for failures, errors, and other non-ideal conditions in both the hardware and the software of a HA system and for monitoring when the open source HA tools detect failures and errors. Further, there is a need for remotely accumulating the monitored system status and then remotely facilitating reconfiguration of the HA system.

[0010] Moreover, typically multiple HA systems are combined in a network to form an enterprise system. Each HA system may service transactional requests for a different store within an enterprise, for example. There is a need for a method, system, and program for remotely accumulating the monitored system status of multiple HA systems within an enterprise, comparing the system status with performance requirements, and tracking hardware and software needs of each HA system within the enterprise.

[0011] Further, when implementing an HA system using an open source operating system framework, it would be advantageous to implement an open source compliant middleware layer to handle transaction requests. In particular, it would be advantageous to implement a Java™ 2 platform, Enterprise Edition (J2EE) compliant middleware stack that is: (1) controlled by open source based cluster management interfacing with a remote enterprise console; and (2) able to monitor and configure multiple HA systems in an enterprise network.

SUMMARY OF THE INVENTION

[0012] The present invention provides improved high availability cluster management and in particular provides for remote cluster management of a high availability system

5 implemented in compliance with an open source framework. Still more particularly, the present invention relates to improved remote monitoring and management of multiple high availability systems in an enterprise network.

[0013] According to one aspect of the present invention, multiple high availability systems are networked in an enterprise and managed overall by a remote enterprise server.

10 Within each high availability system, a cluster management controller monitors a status of a particular component of the high availability system and reacting to adjust the high availability system when the status indicates an error. In addition, with each high availability system, a monitoring controller detects when the cluster management controller reacts to the status of the particular component and detects a condition of a multiple components of the high availability
15 system. The monitoring controller then reports the error and the condition of the components to the remote enterprise server. The remote enterprise server is enabled to manage the high availability system based on the report.

[0014] In particular, the high availability server implement a J2EE compliant middleware stack monitored by open source functions such as a heartbeat monitor and a service
20 monitoring daemon. The heartbeat monitor detects, in particular, the status of particular servers on which the middleware stack resides. The service monitoring daemon detects, in particular the

status of the particular instances of services provided by the middleware stack.

[0015] The remote enterprise server may determine from the report that a configuration change should be made and send a configuration request to the high availability system. The monitoring controller then adjusts the configuration of the high availability system to adjust how
5 the heartbeat monitor or service monitoring daemon will detect and react to errors. Further, other hardware and software components within the high availability system may be reconfigured by the monitoring controller.

[0016] The remote enterprise server preferably stores monitored information about each high availability system in a database. In addition, the enterprise server preferably analyzes the
10 monitored information and determines which high availability systems are not meeting performance requirements. The enterprise server may recommend hardware and software changes and configuration changes. In addition, the enterprise server may display the comparative performance and provide a real-time display of the high availability systems and when errors are detected at each system.

BRIEF DESCRIPTION OF THE DRAWINGS

[0017] The novel features believed characteristic of the invention are set forth in the appended claims. The invention itself however, as well as a preferred mode of use, further
5 objects and advantages thereof, will best be understood by reference to the following detailed description of an illustrative embodiment when read in conjunction with the accompanying drawings, wherein:

[0018] **Figure 1** is a block diagram depicting a server system in which the present
10 method, system, and program may be implemented;

[0019] **Figure 2** is a block diagram depicting the hardware configuration of a high availability cluster for efficient transition of middleware during failover;

15 [0020] **Figure 3** is a block diagram depicting a cluster manager in accordance with the method, system, and program of the present invention;

[0021] **Figure 4** is a block diagram depicting one embodiment of a software configuration of a HA cluster before failover in accordance with the method, system, and
20 program of the present invention;

[0022] **Figure 5** is a block diagram depicting one embodiment of a software configuration of a HA cluster after failover in accordance with the method, system, and program of the present invention;

5 [0023] **Figure 6** is a block diagram depicting one embodiment of an implementation of an independent software vendor application within a J2EE compliant middleware in a HA system;

[0024] **Figure 7** is a high level logic flowchart depicting a process and program for
10 configuring a drbd partition to a J2EE compliant middleware stack in a HA cluster;

[0025] **Figure 8** is a high level logic flowchart depicting a process and program for controlling configuration and failover of a J2EE compliant middleware stack in a HA cluster through a heartbeat controller;

15

[0026] **Figure 9** is a high level logic flowchart depicting a process and program for controlling a mon function for monitoring services provided by a J2EE compliant middleware stack;

20 [0027] **Figure 10** is a block diagram depicting an enterprise network including multiple HA systems running J2EE middleware stacks in accordance with the method, system, and

program of the present invention; and

[0028] **Figure 11** is a high level logic flowchart depicting a process and program for controlling a monitoring controller within a HA cluster manager in accordance with the method,

5 system, and program of the present invention; and

[0029] **Figure 12** is a high level logic flowchart depicting a process and program for remotely controlling a cluster manager of an HA system to reconfigure the HA system; and

10 [0030] **Figure 13** is a high level logic flowchart depicting a process and program for controlling a remote enterprise console for managing multiple HA systems in a cluster.

DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENT

[0031] Referring now to the drawings and in particular to **Figure 1**, there is depicted one embodiment of a system through which the present method, system, and program may be implemented. The present invention may be executed in a variety of systems, including a variety of computing systems, server systems, and enterprise systems.

[0032] Computer system **100** includes a bus **122** or other communication device for communicating information within computer system **100**, and multiple processors **112a-112n** coupled to bus **122** for processing information. Bus **122** preferably includes low-latency and higher latency paths that are connected by bridges and adapters and controlled within computer system **100** by multiple bus controllers.

[0033] Processor **112a-112n** may be a general-purpose processor such as IBM's PowerPC™ processor that, during normal operation, processes data under the control of operating system and application software accessible from a dynamic storage device such as random access memory (RAM) **114** and a static storage device such as Read Only Memory (ROM) **116**. In a preferred embodiment, multiple layers of software contains machine executable instructions that when executed on processors **112a-112n** carry out the operations depicted in the flowcharts of **Figures 7, 8, 9, 11, 12, 13** and others described herein.

Alternatively, the steps of the present invention might be performed by specific hardware components that contain hardwired logic for performing the steps, or by any combination of programmed computer components and custom hardware components.

[0034] The present invention may be provided as a computer program product, included on a machine-readable medium having stored thereon the machine executable instructions used to program computer system **100** to perform a process according to the present invention. The term “machine-readable medium” as used herein includes any medium that participates in providing instructions to processors **112a-112n** or other components of computer system **100** for execution. Such a medium may take many forms including, but not limited to, non-volatile media, volatile media, and transmission media. Common forms of non-volatile media include, for example, a floppy disk, a flexible disk, a hard disk, magnetic tape or any other magnetic medium, a compact disc ROM (CD-ROM) or any other optical medium, punch cards or any other physical medium with patterns of holes, a programmable ROM (PROM), an erasable PROM (EPROM), electrically EPROM (EEPROM), a flash memory, any other memory chip or cartridge, or any other medium from which computer system **100** can read and which is suitable for storing instructions. In the present embodiment, an example of a non-volatile medium is mass storage device **118** which as depicted is an internal component of computer system **100**, but will be understood to also be provided by an external device. Volatile media include dynamic memory such as RAM **114**. Transmission media include coaxial cables, copper wire or fiber optics, including the wires that comprise bus **122**. Transmission media can also take the form of acoustic or light waves, such as those generated during radio frequency or infrared data communications.

[0035] Moreover, the present invention may be downloaded as a computer program product, wherein the program instructions may be transferred from a remote computer such as a

server **140** to requesting computer system **100** by way of data signals embodied in a carrier wave or other propagation medium via one of network links **134a-134n** to a communications interface **132** coupled to bus **122**. Communications interface **132** provides a two-way data communications coupling to multiple network links **134a-134n** that may be connected, for example, to a local area network (LAN), wide area network (WAN). When implemented as a server system, computer system **100** typically includes multiple communication interfaces accessible via multiple peripheral component interconnect (PCI) bus bridges connected to an input/output controller. In this manner, computer system **100** allows connections to multiple network computers.

[0036] In a network environment, computer system **100** communicates with other systems through network **102**. Network **102** may refer to the worldwide collection of networks and gateways that use a particular protocol, such as Transmission Control Protocol (TCP) and Internet Protocol (IP), to communicate with one another. Network **102** uses electrical, electromagnetic, or optical signals that carry digital data streams. The signals through the various networks and the signals on network links **134a-134n** and through communication interface **132**, which carry the digital data to and from computer system **100**, are exemplary forms of carrier waves transporting the information. Although not depicted, computer system **100** may also include multiple peripheral components that facilitate communication.

[0037] When computer system **100** is implemented as a server system in a HA cluster, additional network adapters may be included for supporting local connections with other server systems. In addition, when implemented as a server system in a HA cluster, computer system

100 may be designed as a commodity hardware server, such as the xSeries™ servers from IBM Corp.

[0038] Those of ordinary skill in the art will appreciate that the hardware depicted in **Figure 1** may vary. Furthermore, those of ordinary skill in the art will appreciate that the depicted example is not meant to imply architectural limitations with respect to the present invention.

[0039] With reference now to **Figure 2**, there is depicted a block diagram of the hardware configuration of a high availability cluster for efficient transition of middleware during failover. As illustrated, client systems **202** and **204** are connected to network **102** for transferring requests for services. In the embodiment, client systems **202** and **204** request services from a high availability (HA) system **208** configured for a quick recovery time with minimal committed data loss during failover.

[0040] As illustrated, HA system **208** includes a primary node **210** and a secondary node **220**. As will be further described, primary node **210** and secondary node **220** preferably implement redundant hardware and software that when executed, provide a high availability system. In particular, primary node **210** and secondary node **220** implement redundant middleware which in a preferred embodiment supports J2EE applications, as will be further described. Middleware is the software that develops, integrates, and manages web applications and systems. As will be further described, the middleware enables integration of communication, processes, data and the automation of transaction capacity and systems management.

[0041] In particular, Java™ 2 platform, Enterprise Edition (J2EE) provides a reusable component model for use in building web applications. J2EE defines a standard application model, a standard platform for hosting applications, a compatibility requirement and an operation definition of the J2EE platform. An advantage of this open source model is that multiple
5 developers can implement the J2EE model with additional components and configurations, yet all J2EE applications will run on a J2EE based system.

[0042] Developers at International Business Machines, Corp. (IBM™), have developed software that implements the J2EE model. This software often fills in gaps not specified in the J2EE framework. For example, IBM™, in particular, has developed a middleware stack of J2EE
10 compliant software products that when implemented on a cluster of servers, support J2EE applications. In general, the middleware stack includes a web server, a database server, and a universal Internet application server. Specifically, this stack may include products such as the IBM DB2™ UDB Enterprise Edition, the IBM HTTP Server, and the IBMWebSphere™ Application Server.

15 [0043] In addition, primary node 210 and secondary node 220 implement monitoring and configuration controllers that monitor the failures and errors of J2EE compliant middleware stack and hardware in a HA cluster. As an example of monitoring and configuration controllers, Tivoli™ Monitoring controllers may be implemented that fill in the gaps for monitoring software running in the J2EE framework and facilitate the configuration of systems running the J2EE
20 framework.

[0044] Primary node 210 and secondary node 220 are connected in a simple, reliable

manner that enables each node to quickly check the heartbeat of the other node. In the embodiment, this connection is enabled by a cross-over cable **218** connected between network adapters at each node. In particular, cross-over cable **218** preferably enables an Ethernet connection for transferring heartbeat data. Alternatively, heartbeat data may also be transferred
5 across the public IP connection via network **102** in event that cross-over cable **218** fails. It will be understood that other hardware may be implemented for providing the heartbeat communication channel between primary node **210** and secondary node **220** and that in addition to a network based connection, a serial connection may be implemented.

[0045] In particular, when a heartbeat signal is sent between primary node **210** and
10 secondary node **220** over cross-over cable **218**, if the heartbeat fails, then secondary node **220** will take over the services provided by primary node **210** before the failure. As will be further described, however, according to an advantage of the present invention, middleware components may further analyze the heartbeat failure and provide additional information about the failure before secondary node **220** takes over for the services provided by primary node **210**. Further, as
15 will be further described, both Linux based and non-Linux based heartbeats may be monitored via cross-over cable **218**.

[0046] Primary node **210** and secondary node **220** access data storage systems **214** and
20 **224**. Advantageously, a data replicator, herein depicted as a drbd partition **230**, including a partition of each of data storage systems **214** and **224**, for replicating data accessible by primary node **210** and secondary node **220** without requiring a storage device that is actually physically shared between primary node **210** and secondary node **220**. According to an advantage of the

present invention, drbd is configured to run on the partition to facilitate the transfer of data during failover from primary node **210** to secondary node **220**. It will be understood that while the invention is described with respect to a drbd partition managed by drbd scripts, other distributed data replication systems may be implemented.

5 **[0047]** Uninterrupted power supply (UPS) **212** and UPS **222** each provide an independent power supply to primary node **210** and secondary node **220**, respectively. Preferably, a connection is also established between UPS **212** and secondary node **220** and UPS **222** and primary node **210**. In one embodiment, a serial cable **216** is provided from primary node **210** to UPS **222** and a serial cable **226** is provided from secondary node **220** to UPS **212**. It will
10 be understood, however, that other types of connection hardware may be implemented.

[0048] According to an advantage of the present invention, when a failure is detected in primary node **210**, secondary node **220** begins receiving the requests previously directed to primary node **210** after failover. Because only a portion of the hardware, software, or network running on primary node **210** may fail, the only way to ensure that primary node **210** does not try
15 to update data after the failover is to turn off UPS **212**. Advantageously, as will be further described, when the failover to standby node **220** is detected, STONITH, described in more detail herein, is implemented by the cluster manager to direct a command from standby node **220** to UPS **212** to turn off the power supply.

20 **[0049]** With reference now to **Figure 3**, there is depicted a block diagram of a cluster manager in accordance with the method, system, and program of the present invention. As

illustrated, a cluster manager **322** includes multiple components utilized to implement an efficient failover including a heartbeat tool **402**, drbd scripts **404**, mon **406**, and a stonith function **408**. It will be understood that other components may be included in a cluster manager to manage other aspects of the cluster. Further, it will be understood that additional components may be included in cluster manager **322** to manage failover.

[0050] Heartbeat tool **402** preferably includes the heartbeat package for Linux, configured for managing failover within a HA cluster with a J2EE compliant middleware stack. In particular, Heartbeat tool **402** generally works by sending a “heartbeat” request between two nodes in a cluster. As described in **Figure 2**, the heartbeat request may be sent through cross-over cable between network adapters at each node. When applied to a J2EE compliant middleware stack running on clusters of server systems, heartbeat requests sent by heartbeat tool **402** are distributed about the different layers of the stack.

[0051] If the heartbeat request fails to be returned, then the secondary node can assume that the primary node failed and take over IP, data, and services that were running on the primary node. When the secondary node takes over the IP, data, and services that were running on the primary node, heartbeat tool **402** startups components of the secondary node that are waiting in standby mode, assigns IP addresses to components of the secondary node, and performs other failover tasks.

[0052] Drbd **404** is a kernel module with associated scripts that that manage data in a HA cluster for improved switching of data during failover. This is performed by mirroring a block device managed by drbd **404**. Drbd is a script that loads the drbd module and configures

with the IP addresses of the relevant systems in the HA cluster and the shared storage device.

[0053] When applied to a J2EE compliant middleware stack, the drbd managed block device provides storage on which the middleware stack can run. Initially, the cluster is configured and the drbd partition is mounted so that only the primary node can read or write from the drbd managed block device. When a failover occurs, the datadisk script of drbd **404** is run by heartbeat tool **402** to mount the drbd partition so that only the secondary node can read/write from the drbd managed block device.

[0054] Mon **406** is a service monitoring daemon that periodically runs monitoring scripts that monitor critical system services within the J2EE compliant middleware stack. If a service is found to have failed or terminated abnormally, mon **406** restarts the service to ensure that all components of the middleware stack remain running within the primary service. Abnormal termination may occur, for example, from programming errors or catastrophic operating system events such as temporary critical resource constraints with RAM. In particular, when mon restarts a service, it restarts a new instance of the service with a process identifier (PID) different from the dead service, but the same virtual IP address.

[0055] Stonith **406** is a function called by heartbeat tool **402** to ensure data integrity during failover. In particular, stonith **406** includes the configuration of the serial cables to UPS **212** and **222**, as depicted in **Figure 2**. When heartbeat tool **402** calls stonith **406**, the call designates the node to be shutdown. Stonith sends a signal to turn off the power of the requested UPS.

[0056] Monitoring and configuration controller **410** includes multiple monitoring

controllers which are specified for monitoring the status of hardware and software within the HA clusters. According to an advantage of the invention, status information about the multiple hardware and software components of HA clusters is forwarded to a remote centralized enterprise console. Preferably, monitoring and configuration controller **410** supplements the

5 Java™ Management Extensions (JMX) to monitor the hardware and software components of the HA clusters, to detect bottlenecks and potential problems, and to automatically recover the cluster from critical situations. In one embodiment, the monitoring controllers are enabled by Tivoli™ Monitoring which forwards monitored information to a Tivoli™ Enterprise Console (TEC).

10 [0057] In particular, while heartbeat tool **402** and mon **406** monitor the status of specific components and specific instances of services within the nodes, monitoring and configuration controller **410** detects the conditions monitored by these tools and detects the overall status of the system when heartbeat tool **402** is triggered to initiate failover or mon **406** is triggered to restart a server. Thus, monitoring and configuration controller **410** supplements the

15 open source tools by compiling the status of multiple components of the nodes when failures, errors, and non-ideal conditions occur.

[0058] According to one advantage of the invention, the remote centralized monitoring console can use the information gathered to determine configuration changes. In particular, according to an advantage of the invention, the monitoring controllers of monitoring and

20 configuration controller **410** are each configured to monitor each hardware component in the HA cluster and each of the layers of the J2EE compliant middleware stack. Thus, based on

monitored information about the hardware and middleware layers, the console can determine which middleware layers need more memory for caching requests, need more threads for handling requests, or need to be reconfigured in some other manner. The console can send configuration changes to the configuration controllers of monitoring and configuration controller 5 410, which then adjust the configuration of the HA clusters. In one embodiment, the configuration controller is a Tivoli™ Configuration Manager which manages the configuration characteristics of the HA clusters.

[0059] According to another advantage of the invention, in an enterprise system, the console use the information gathered to determine which HA clusters need hardware and 10 software upgrades. For example, for the monitored information, the console can determine which stores have hardware which seems to be failing and needs to be replaced, which stores have hardware which has reached capacity and needs to be upgraded, and which stores have software that is failing or not running reliably.

[0060] According to yet another advantage of the invention, monitoring and 15 configuration controller 410 interacts with the other monitoring components within cluster manager 322 to gather the status information that is sent to the console. For example, when mon 406 detects a failure of any of the monitored services, monitoring and configuration controller 410 sends a notification to the remote centralized monitoring console so that a bigger picture of failures in the system can be compiled. Further, when heartbeat tool 402 initiates a failover of 20 one node of the system to another node, monitoring and configuration controller 410 sends a notification to the remote centralized monitoring console so that node failure statistics can be

gathered.[0060] With reference now to **Figure 4**, there is depicted a block diagram of one embodiment of a software configuration of a HA cluster before failover in accordance with the method, system, and program of the present invention. As depicted, primary node **210** and secondary node **220** represent clusters of server systems, each assigned to an IP address.

5 [0061] According to an advantage of the present invention, cluster manager **322** runs on primary node **210** and secondary node **220** to monitor for failures, restart services, and control failover when a failure is detected. As illustrated, cluster manager **322** sets up drbd partition **230** that is located on storage shared between primary node **210** and secondary node **220**.

 [0062] Primary node **210** includes all active components of the middleware stack: a
10 load balancer **312**, HTTP servers **314**, web application servers (WAS) **316**, messaging controllers **318**, and a database server **320**. Secondary node **220** includes active HTTP servers **334** and WASs **336**, however, load balancer **332**, messaging controllers **338**, and database **340** are in standby mode.

 [0063] Load balancers **312** and **332** preferably balance the load of requests between
15 HTTP and WAS servers, which may also be clustered. Preferably, load balancers **312** and **314** perform intelligent load balancing by using server availability, capability, workload, and other criteria. According to one embodiment, load balancers **312** and **332** may be implemented through the IBM WebSphere™ Edge Server.

 [0064] As illustrated, load balancers **312** and **332** may implement a heartbeat
20 independent of the Linux based heartbeat. Alternatively, the Linux based heartbeat monitoring **332** and **342** may monitor the status of load balancers **312** and **332**.

[0065] HTTP servers **314** and **334** may include clusters of servers designed to receive HTTP requests and distribute HTTP requests among WAS **316** and **336**, respectively. In addition, HTTP servers **314** and **334** are enabled to call enablers, such as servlet containers and Enterprise Java™ Bean (EJB) containers, when other requests, such as requests for servlets and EJBs, are received. According to one embodiment, HTTP servers **314** and **334** may be implemented through an HTTP server bundled with IBM's WebSphere™, and in particular WebSphere™ v. 5.0. WebSphere™ 5.0 is advantageous because multiple copies of the WebSphere™ components can be controlled from one location. Thus, configuration changes can be made in one place that effects multiple instances of the software components located on multiple server systems.

[0066] According to an advantage of the present invention, HTTP servers **314** and **334** are run in an active/active configuration where the heartbeat tool of cluster manager **322** activates HTTP server after primary node is up and running. By running HTTP servers **314** and **334** in an active/active configuration, the request load can be split across the two (or more) servers to increase the speed at which client requests are handled. In addition, by running HTTP servers **314** and **334** in an active/active configuration, then startup time on failover is reduced.

[0067] WAS **316** and **336** preferably include clusters of servers enabled to support web applications providing mission-critical services to customers, and in particular these servers are enabled to support J2EE applications. According to one embodiment, WAS **316** and **336** are WebSphere™ Application Servers supported by IBM's Websphere™ 5.0 that host the servlets, EJBs, and other J2EE components necessary for supporting a J2EE application and services.

[0068] WAS 316 interacts with messaging controller 318 and database server 320 to provide application server functionality integrated with messaging control and databases.

According to an advantage of the present invention, WAS 316 and WAS 336 are run in an active/active configuration. In particular, when initializing the systems, once messaging controller

5 318 and database server 320 are available, the heartbeat tool of cluster manager 322 launches WAS 336 to create the active/active configuration. By running an active-active configuration, the request load can be split across multiple clusters of systems to increase the speed at which client requests are handled. In addition, by running an active/active configuration, then startup time on failover is reduced.

10 [0069] Messaging controllers 318 and 338 include a controller for listening for asynchronous requests and storing those requests in a local queue to provide a queue to communicate with J2EE based systems. Messaging controller 318 and 338 may implement IBM MQSeries™, IBM WebSphere™ MQ, or other message controllers that supplement the Java™ Messaging Service (JMS).

15 [0070] According to an advantage of the present invention, messaging controllers 318 and 338 are run in an active/standby configuration where the drbd of cluster manager 322 manages the persistent resources in the messaging queue in drbd partition 230 and the heartbeat tool of cluster manager 322 controls the startup of messaging controller 338 in a failover.

[0071] Database servers 320 and 340 provide control for persistent storage. Database
20 servers 320 and 340 may be implemented through a database control system such as IBM DB2 UDB Enterprise Edition or other relational database management systems.

[0072] According to an advantage of the present invention, database servers 320 and 340 are run in an active/standby configuration where the drbd of cluster manager 322 manages the persistent resources in the database in drbd partition 230 and the heartbeat tool of cluster manager 322 controls the startup of database server 340 in a failover.

5 [0073] For messaging controllers 318 and 338 and database servers 320 and 340 to run in active/standby configuration and quickly failover with minimal data loss, messaging controller 318 and database server 320 are configured to point to the location where drbd partition 320 is mounted as the root for storage of the queue and database. In addition, cluster manager 322 configures drbd and the heartbeat tool with the virtual IP address of messaging controller 318 and
10 database server 320.

[0074] Further, according to an advantage of the present invention, the mon function of cluster manager 322 periodically runs monitoring scripts that monitor critical system services, such as the services provided by messaging controller 318 and database server 320. If a service is found to have failed or terminated abnormally, mon restarts the service to ensure that all
15 components of the middleware stack remain running within the primary service.

[0075] It is important to note that the method of configuring each level of middleware to achieve efficient failover and controlling each level of middleware through cluster manager 322 may be applied to other types of middleware. Thus, as the functions available from a middleware software stack that is J2EE compatible continue to expand, each middleware
20 component can be configured either in an active/active or active/standby configuration, monitored by cluster manager 322, and controlled during failover.

[0076] With reference now to **Figure 5**, there is depicted a block diagram of one embodiment of a software configuration of a HA cluster after failover in accordance with the method, system, and program of the present invention. As depicted, after failover, primary node 210 is marked as a failed node. Secondary node 220 takes over as the all active node.

[0077] When a failure is detected and secondary node 220 designates primary node 210 as “dead”, hardware and software issues are present. In particular, primary node 210 may not respond to a heartbeat request within the required time, however, would be operational shortly thereafter. To avoid a case where primary node 210 and secondary node 220 are both operational, as previously described, the heartbeat tool of cluster manager 322 will call STONITH to turn off the UPS to primary node 210. By implementing an inexpensive UPS, controllable by STONITH, data integrity can be achieved and the “split brain” issue of HA that can arise when the primary node is not really dead is avoided.

[0078] Next, during failover, the load balancer heartbeat manages startup of load balancer 332. When activated, the heartbeat tool of cluster manager 322 assigns the virtual IP1 address of primary node 210 to load balancer 332.. Accordingly, requests to the virtual IP address are redirected to load balancer 332 so that no change in the IP address of the load balancing cluster occurs.

[0079] During failover, since HTTP server 334 and WAS 336 are already active, the heartbeat tool of cluster manager 322 does not need to startup these components. However, since messaging controller 338 and database server 340 are in standby, the heartbeat tool of cluster

manager **322** needs to manage the failover of these layers. First, the heartbeat tool will takeover the virtual IP2 address. Next, the heartbeat tool will start the datadisk service of the drbd to configure and mount the drbd mirrored partition. Finally, the heartbeat tool will startup messaging controller **338** and database server **340** configured to virtual IP2 address and with the message queue and database instances launching on mirrored drbd partition **230**. Alternatively, although not depicted, database server **340** may be in active mode, rather than standby, because the virtual IP2 address is only available to one node at a time. Since database server **340** does not attempt to touch the data on drbd partition **230** until a request arrives, at failover, database server **340** is configured to the virtual IP2 address and mirrored drbd partition **230** is accessible before a request arrives. In contrast, some layers, such as messaging controller **338** load data directly at startup and thus will crash if started up on secondary node **220** before failover because data on drbd partition **230** is not available to secondary node **220** before failover. .

[0080] Referring now to **Figure 6**, there is depicted a block diagram of one example of an implementation of an independent software vendor (ISV) application within a J2EE compliant middleware in a HA system. As illustrated, an active WAS **602**, active IBM MQSeries™ server **610**, and active IBM DB2 server **614** illustrate a portion of the primary node of a J2EE compliant middleware stack interfacing with a drbd partition **630**. As illustrated at reference numeral **620**, an item sale or transaction completion is received at an active WebSphere™ Application Server **602**. An ISV may program a servlet or EJB to handle a particular type of incoming request. For example, as depicted at reference numeral **620**, a lookup servlet **604** is an

ISV web application that handles price lookups (PLU) to check the price of items as they are scanned at a cash register. Lookup servlet **602** then posts a request for the retain transaction to be completed asynchronously by another component, such as transaction servlet **608** or another servlet or EJB. First, however, as depicted at reference numeral **622**, the information is transferred to MQ listener **612** and placed on MQ queue **632** to free lookup servlet **604** to receive the next incoming request and to ensure that the transaction will be recorded exactly once, in order, via MQ queue **632**. Next, as depicted at reference numeral **624**, MDB **606** is then called to take the transaction off MQ queue **632** and as depicted at reference numeral **626**, to feed the transaction to transaction servlet **626**. Transaction servlet **626** ultimately processes the PLU and, as depicted at reference numeral **628**, commits the result to IBM DB2 controller **616** for storage in DB2 **634**.

[0081] In particular, **Figure 6** illustrates the advantages of the J2EE compliant middleware stack in a HA system during failover because the stack ensures that each transaction will be recorded exactly once, even if failover occurs after a request has already begun to transition between the layers of the stack. In addition, **Figure 6** illustrates the advantages of the J2EE compliant middleware stack in a HA system during failover because active layers MQSeries™ server **610** and DB2 server **614** interface with drbd partition **630** that is only accessible to the primary node, but is quickly remounted for access by the secondary node during failover.

[0082] With reference now to **Figure 7**, there is depicted a high level logic flowchart of

a process and program for configuring a drbd partition to a J2EE compliant middleware stack in a HA cluster. As depicted, the process starts at block 700 and thereafter proceeds to block 702. Block 702 depicts configuring and mounting the drbd partition. Next, block 704 depicts activating the message queue and database on the drbd partition. Thereafter, block 706 illustrates recording the virtual IP address of the messaging server and database server accessing the drbd partition for efficient transfer of access to the drbd partition during failover, and the process ends.

[0083] Referring now to **Figure 8**, there is depicted a high level logic flowchart of a process and program for controlling configuration and failover of a J2EE compliant middleware stack in a HA cluster through a heartbeat controller. As depicted, the process starts at block 800 and thereafter proceeds to block 802. Block 802 illustrates activating the middleware layers of the primary node. Thereafter, block 804 depicts activating the HTTP server and the WAS middleware layers of the secondary node. In addition, other middleware layers that are designated to run in an active-active configuration are activated. Thereafter, block 806 depicts periodically initiating a heartbeat request from the secondary node to the primary node. Block 808 depicts a determination whether a heartbeat return is detected by the secondary node. If a heartbeat return is detected, then the process returns to block 806. If a heartbeat return is not detected, then the process passes to block 810.

[0084] Block 810 depicts calling STONITH to turn off the power supply of the primary node. Next, block 812 depicts taking over the virtual IP addresses from the primary node to assign to the redundant component in the secondary node. Thereafter, block 814 depicts calling

the datadisk script to remount the drbd partition for access by the secondary node and the process ends. Then, block 816 depicts activating the standby middleware layers on the secondary node and launch data on the drbd partition. It will be understood that additional steps may be performed by the heartbeat tool and other cluster management services during failover.

5

[0085] With reference now to **Figure 9**, there is depicted a high level logic flowchart of a process and program for controlling a mon function for monitoring services provided by a J2EE compliant middleware stack. As depicted, the process starts at block 900 and thereafter proceeds to block 902. Block 902 depicts configuring a schedule for monitoring services provided by the
10 middleware. Next, block 904 depicts a determination whether a scheduled monitoring time is triggered. If a scheduled monitoring time is not triggered, then the process iterates at block 904. If a scheduled monitoring time is triggered, then the process passes to block 906. Block 906 depicts monitoring the status of the scheduled service. Thereafter, block 908 depicts a determination whether the service is detected as dead or failed in some manner. If the service is
15 not detected as dead, then the process ends. If the service is detected as dead, then the process passes to block 910. Block 910 depicts restarting the same service with a new PID, and the process ends.

[0086] Referring now to **Figure 10**, there is depicted a block diagram of an enterprise
20 network including multiple HA systems running J2EE middleware stacks and managed by a remote enterprise console in accordance with the method, system, and program of the present

invention. As depicted, a HA system **1202** and a HA system **1204** are communicative connected to a remote enterprise console **1210** that monitors and remotely controls HA systems **1202** and **1204** via network **102**. It will be understood that multiple HA systems may be monitored and controlled by a single or multiple remote central consoles.

5 **[0087]** According to an advantage of the invention, each of HA systems **1202** and **1204** may handle retail transactions and other mission critical operations. According to one embodiment, each of HA systems **1202** and **1204** enable high availability through redundant J2EE compliant middleware stacks that enable J2EE applications, such as the middleware stacks illustrated in **Figures 4** and **5**. In particular, each of HA systems **1202** and **1204** includes a
10 cluster manager running monitoring and configuration controllers **410**, as depicted in **Figure 3**.

[0088] Advantageously, when errors, failures, or non-ideal conditions occur at any of HA systems **1202** and **1204**, monitoring and configuration controllers **410** detects the condition of the system at the time of the error, failure or other non-ideal condition and then compiles the information to make a report to remote enterprise console **1210**. According to an advantage of
15 the invention, if the heartbeat monitor or mon functions detect a failure or error, then monitoring and configuration controllers **410** are triggered to detect the failure or error and determine the system conditions at the time of the failure or error.

[0089] Remote enterprise console **1210** preferably stores monitored information in a database. Next, remote enterprise console **1210** preferably includes a first controller for
20 analyzing the error and failure information received from HA systems **1202** and **1204** and potentially returning configuration changes to the HA systems to attempt to prevent and improve

the efficiency of failovers. In addition, remote enterprise console **1210** may include a second controller that compares the failures, errors and other information received from multiple HA systems to determine which systems need repairs and upgrades and which systems are not meeting performance requirements. Remote enterprise console **1210** may gather and control display of performance statistics for HA systems **1202** and **1204**.

[0090] Referring now to **Figure 11**, there is depicted a high level logic flowchart of a process and program for controlling a monitoring controller within a HA cluster manager in accordance with the method, system, and program of the present invention. As depicted, the process starts at block **1000** and thereafter proceeds to block **1002**. Block **1002** depicts a determination whether a failure or error is detected from the heartbeat monitor, mon, or other monitoring controller monitoring a middleware stack in a HA system. If no failure or error is detected, then the process iterates at block **1002**. If a failure or error is detected, then the process passes to block **1004**. Block **1004** depicts gathering and analyzing available system information at the time of the failure or error. Next, block **1006** depicts sending the failure or error and available system information to a remote central console monitoring the HA system, and the process ends.

[0091] With reference now to **Figure 12**, there is depicted a high level logic flowchart of a process and program for remotely controlling a cluster manager of an HA system to reconfigure the HA system. As illustrated, the process starts at block **1100** and thereafter

proceeds to block **1102**. Block **1102** depicts a determination of whether a configuration request is received from a remote enterprise console to reconfigure the HA system running a middleware stack. If the request is not received, then the process iterates at block **1102**. If the request is received, then the process passes to block **1104**. Block **1104** depicts calling the heartbeat monitor to reconfigure the HA system failover settings, and the process ends. In addition, other controllers within the cluster manager of a HA system may be called to adjust other software and hardware configurations of the HA system.

[0092] Referring now to **Figure 13**, there is depicted a high level logic flowchart of a process and program for controlling a remote enterprise console for managing multiple HA systems in a cluster. As depicted, the process starts at block **1300** and thereafter proceeds to block **1302**. Block **1302** depicts a determination whether monitored information is received from an HA system. If monitored information is not received, then the process iterates at block **1302**. If monitored information is received, then the process passes to block **1304**. In particular, the remote enterprise console may periodically send request to each of the HA system for monitored information and each HA system may also automatically send monitored information.

[0093] Block **1304** depicts adding the monitored information to an enterprise database in which monitored information from multiple HA systems is stored. Next, block **1306** depicts requesting reconfiguration of the HA system if the monitored information triggers reconfiguration. In particular, a remote enterprise console may include predetermined configurations to be requested when particular types of errors are detected in monitored

information. Alternatively, a system administrator may recommend the type of configuration for a particular type of error. Thereafter, block 1308 depicts recalculating the performance statistics for the HA system based on the monitored information. In particular, calculation of performance statistics may only be triggered for certain types of monitored errors or fluctuations. Next, block

5 1312 depicts comparing the performance of this HA system with the performance of the other HA systems in the enterprise network and the performance requirements set for the enterprise network. Then, block 1314 depicts displaying the comparative performance results in charts and graphs. For example, a chart may depict graphical representations of the locations of the HA systems and provide a graphical indicator of which systems have failed and provide graphical
10 indicators to show the performance of each HA system relative to the other HA systems. Further, real-time performance of each system and any errors reported may be displayed. Next, block 1316 depicts recommending corrective action for HA system weaknesses, and the process ends. For example, the recommendations may indicate which HA systems need to be replaced, which HA systems need to be upgraded, and which HA systems need software upgrades or fine tuning.
15 It will be understood that the processes depicted in **Figure 13** are examples of types of processes that can be performed on the monitored information received from multiple high availability servers and that other similar analysis and outputs can be performed without departing from the scope of the invention.

20 [0094] While the invention has been particularly shown and described with reference to a preferred embodiment, it will be understood by those skilled in the art that various changes in

form and detail may be made therein without departing from the spirit and scope of the invention.